

***A Flexible Run-time Support for Distributed
Dependable Hard Real-time Applications***

Emmanuelle Anceaume, Gilbert Cabillic

Pascal Chevochot and Isabelle Puaut

N° 3564

Novembre 1998

THÈME 1

 ***apport
de recherche***

A Flexible Run-time Support for Distributed Dependable Hard Real-time Applications

Emmanuelle Anceaume, Gilbert Cabillic

Pascal Chevochot and Isabelle Puaut

Thème 1 — Réseaux et systèmes

Projet Solidor

Rapport de recherche n3564 — Novembre 1998 — 27 pages

Abstract: Typically, most distributed, dependable, real-time systems designed in the past can only meet the particular requirements of the application domain to which they were targeted. This approach led to specific, non-flexible, dedicated and non-reusable solutions, often based on specialized hardware. This report presents an alternative approach where a flexible run-time support for distributed dependable hard real-time applications is built on top of off-the-shelf hardware. This support, called HADES has been designed by considering three fundamental and complementary aspects: *real-time*, to support applications that exhibit hard timing constraints; *fault-tolerance*, to provide a high degree of reliability through transparent fault tolerant software components; and *flexibility*, to allow the modifications of components of the run-time support without having to rewrite it entirely, and to support a large range of application domains, real-time kernels and hardware.

Key-words: Run-time support, real-time, availability, distribution, flexibility.

This work is partially supported by the French Department of Defense (DGA/DSP), #96.34.106.00.470.75.65.

Unité de recherche INRIA Rennes

IRISA, Campus universitaire de Beaulieu, 35042 RENNES Cedex (France)

Téléphone : 02 99 84 71 00 - International : +33 2 99 84 71 00

Télécopie : 02 99 84 71 71 - International : +33 2 99 84 71 71

(Résumé : *tsvp*)

Support d'exécution flexible pour des applications temps-réel distribuées hautement disponibles

Résumé : La plupart des systèmes temps-réel fiables développés dans le passé furent dédiés à un domaine d'application et une politique d'ordonnancement donnés. Cette approche a souvent entraîné le développement de solutions lourdes, rigides et reposant fortement sur le développement de matériel spécifique. Ce rapport présente une approche alternative dans laquelle un support d'exécution flexible destiné à développer et exécuter des applications temps-réel distribuées et hautement disponibles est construit au-dessus de matériels *sur étagère*. Ce support, appelé HADES, a été conçu en considérant trois aspects fondamentaux et complémentaires : *temps-réel*, pour exécuter des applications exhibant des contraintes temporelles fortes ; *tolérance aux fautes*, pour fournir un haut degré de fiabilité à travers des mécanismes de tolérance aux fautes transparent à l'applicatif ; et *flexible*, pour permettre d'une part, l'évolution du support d'exécution sans entraîner sa réécriture complète, et d'autre part sa capacité à supporter un large éventail de noyaux temps-réel, de matériels et de domaines d'application.

Mots-clé : Support d'exécution, temps-réel, distribué, disponibilité, flexibilité.

1 Introduction

Dependable hard real-time applications in the last thirty years have been spread out in space rockets and flight control systems, nuclear power plants, stock exchange, medical and automotive equipments. Designing systems supporting such demanding applications is difficult because they must reach a triple goal. First, these systems have to be *highly reliable* to meet the maximum acceptable probability of failure ranging from 10^{-4} to 10^{-10} failures per hour [RSL95, KWFT88]. Second, they have to ensure *strict timeliness* requirements to guarantee response times typically ranging from 1 to 100ms [MRS⁺90, RSL95]. Third, they must enforce *data consistency* to handle concurrent executions of multiple tasks, especially in a distributed environment.

Reaching this triple goal is complex and led so far to solutions that met only the particular requirements of an application domain. By particular requirements we refer to the constraints imposed by the applications in terms of resources consumption, specific scheduling needs, or fault-tolerance. Existing solutions are therefore dedicated to a single application domain, and often based on specialized (and costly) hardware, making the implemented software seldom reusable in a different context [RSL95, KWFT88].

The HADES project¹ [ACCP98] developed at IRISA addresses these problems and provides an environment for the development and the execution of distributed dependable hard real-time applications.

Application development in HADES relies on the use of an original task model and a set of off-line tools. The HADES task model defines the way application tasks are structured and the constraints that must be handled during their execution. The off-line tools encompass a worst-case execution time analyzer, feasibility tests associated with a panel of scheduling policies, an automatic replication tool, and a tool that determines the right amount of memory to be allocated for the application execution. The replication tool [CPC98] automatically replicates application tasks according to the requirements of the application designer (e.g.,

¹HADES stands for Highly Available Distributed Embedded System.

tasks or the portions of tasks to be replicated, replication strategy to be used such as active, passive or semi-active replication).

Application execution is managed by the HADES run-time support. This run-time support is built as a middleware software layer running on top of off-the-shelf real-time kernels (it has been ported on Chorus [CS96] and emulated on Solaris). It consists in a set of mandatory services for the execution of a large panel of dependable distributed real-time tasks. Examples of such services are scheduling policies, communication primitives, group membership management, distributed execution control, and clock synchronization. All these services exhibit timeliness and dependable properties.

The HADES run-time support has been designed to guarantee three fundamental and complementary aspects:

1. *Real-time*: support of applications that exhibit hard timing constraints;
2. *Fault-tolerance*: provision of a high degree of reliability through transparent fault tolerant software components;
3. *Flexibility*, that is:
 - (a) *Extensibility*: ability to change parts of the run-time support without having to rewrite it entirely.
 - (b) *Adaptability*: support for a large range of application domains, real-time kernels and hardware.

The achievement of the real-time aspect mainly relies on an accurate estimation of the worst-case execution times (WCETs) for all the activities to be executed in the system (application tasks, run-time support tasks, real-time kernel). The computation of WCETs of the application and run-time support tasks is eased by the HADES task model (see Section 2); in particular, it requires that all synchronization constraints (e.g., precedence between tasks, resource sharing) be specified statically. Concerning the core of the run-time support, it has been kept small enough to be easily analyzed. Lastly, the predictability of the real-time kernel has been verified through the analysis of the Chorus kernel source code.

The fault tolerance aspect of the run-time support relies on the automatic task replication tool, on the fault-detection and exception handling mechanisms offered by the run-time support. For space consideration, this aspect is not detailed in this paper, but the interested reader is invited to read [CPC98].

To achieve flexibility, we have specified in HADES a set of services with a well-defined interface, that can evolve without having to rewrite the run-time support (*extensibility* property). These services are relevant to task scheduling, communication, group membership, execution management of distributed tasks, and clock synchronization. Any implementation of any service can be replaced by a service with different properties as far as it conforms to the service interface. This allowed us to develop a range of service implementations adapted to the specific constraints of the applications and the hardware (*adaptability* property).

The remainder of the paper is devoted to the description of the HADES run-time support. It concentrates on the main structuring principles that were exploited to obtain a run-time support exhibiting a timely behavior, while maintaining its flexibility. It is organized as follows. Section 2 briefly presents the HADES task model. The HADES run-time support is described in Section 3. This description encompasses an overview of the services interface illustrating the extensibility of HADES (Section 3.2). Section 3.3 points out the completeness of this interface to demonstrate the adaptability of HADES. Finally, the HADES run-time support is compared with related work in Section 4.

2 The HADES task model

The HADES task model describes every task as a directed acyclic graph (see Section 2.1). A set of timing, synchronization, distribution and fault-tolerance attributes can be specified for each task (see Section 2.2). The fault model assumed when designing these tasks is given in Section 2.3. The term *task* covers both the application tasks, denoted by *App_tasks*, and the tasks implementing run-time support services, denoted by *RTS_tasks*. In the remainder of this section, we give only a brief description of the task model. The reader is invited to read [ACCP98] for more details.

2.1 Structure of tasks

Every task is described as a direct acyclic graph, called *HEUG*, for *HADES Elementary Unit Graph*. A node (also called *EU* for *Elementary Unit*) can either be a *Code Elementary Unit* (*Code_EU*), i.e., a sequence of code without synchronization, or a *System Call Elementary Unit* (*Sys_EU*). Tasks can use an *invocation Sys_EU* to request the synchronous or asynchronous execution of other tasks, or to interact with the core of the run-time support. An edge is a precedence constraint between nodes. An edge from an *Elementary Unit* eu_i to an *Elementary Unit* eu_j , noted $eu_i \rightarrow eu_j$, states that eu_j can start executing only after eu_i has finished its execution.

2.2 Task attributes

Task attributes can be classified into synchronization attributes, timing attributes, distribution attributes and fault-tolerance attributes.

Synchronization attributes: The *Code_EUs* can be synchronized by using either condition variables (i.e., boolean variables that can be cleared, set and waited for by *Code_EUs*) or shared resources expressing consumer-producer synchronization schemes between *Code_EUs*. A *resource* abstracts any hardware or software component required to execute a *Code_EU*. A resource can be associated with persistent data, i.e., data that survives beyond the execution of a *Code_EU*. A *Code_EU* or a group of *Code_EUs* must specify which resources to acquire before beginning its execution, together with the requested access modes (*shared* or *exclusive*).

Timing attributes: Timing attributes express temporal properties at the task level and at the elementary unit level. They can be used off-line, for instance by static feasibility tests, and/or on-line for the purpose of execution monitoring, or dynamic scheduling policies. At the *HEUG* level, attributes express either the expected arrival law (i.e., periodic, sporadic, aperiodic) or the relative deadline (i.e., the date at which the task execution must be completed, relative to the task activation date). The scheduling service to be used to

compute the execution order of *App_task EUs* can be chosen by the application designer (see Section 3). At the *EU* level, the application designer can specify the *EU* deadline, earliest start time, latest start time, and execution priority.

Distribution attributes: A task may execute on a set of sites. A site has a processor, a memory, and a set of local hardware devices (peripheral devices, sensors, actuators). Distribution attributes specify for each *Code_EU* and *Sys_EU* the site on which the corresponding *EU* executes. A precedence constraint $eu_i \rightarrow eu_j$ is said to be *local* (resp. *remote*) if eu_i and eu_j are assigned to the same site (resp. different sites). The application designer can specify which execution service is in charge of managing the execution of the *App_tasks* (see Section 3).

Fault-tolerance attributes: Fault-tolerance attributes specify for each *HEUG* or portions of them (i.e., *HEUG* subgraphs) which fault-tolerance strategy is to be applied (active, passive, semi-active, temporal replication) and the requested replication degree. These attributes are used by the HADES off-line replication tool [CPC98]. The application designer can also specify for each *Code_EU* a handler that catches an exception whenever it is raised. An HADES exception is defined as a synchronous event triggered by the run-time support to notify the *Code_EU* of errors related to its execution (e.g. processor exception, attempt to execute an operation on a failed resource). This enables the designer to define specific fault-tolerance strategies.

2.3 Computational fault model

HADES has been designed to support omissions, value and timing faults. We have defined mechanisms to simulate *fail-silent sites*, i.e., sites for which omission, value and timing faults affecting its processor, memory or some other hardware device causes the site to stop before propagating the erroneous state to another site. Mechanisms (see [CPC98]) have been defined at the task construction level (e.g., temporal replication), and at the run-time support level (e.g., fault-detection) to approximate the fail-silent behaviour. The

communication network is subject to omission failures, so that messages might get lost. To tolerate these failures, we assume that the number of successive omissions is bounded.

3 Description of the HADES run-time support

A key concern when designing the HADES run-time support was to guarantee its timeliness while retaining its flexibility. This section gives an overview of the structuring principles that were followed during HADES construction.

3.1 Overall functionalities of the run-time support

The HADES run-time support executes on every site of the distributed system as a middleware software layer based on a Commercial-Off-The-Shelf (COTS) real-time kernel. Its purpose is to provide a range of flexible functionalities to execute dependable distributed real-time applications. The HADES run-time support has been divided into three software layers, namely the *dispatcher*, the *kernel adaptation layer* and the *services layer*. A non-exhaustive view of the functionalities provided by the run-time support is depicted in Figure 1.

Dispatcher: The dispatcher is the core of the run-time support. It provides the very basic functionalities needed to execute a task, namely *task eligibility management* and *task monitoring*. It interacts both with the services layer and the adaption layer.

Concerning eligibility management, the dispatcher maintains a priority-ordered set of run-queues. *App_tasks* are inserted in a certain order in the appropriate run-queues by the *scheduling service* (see Section 3.2). In contrast, *RTS_tasks* are inserted by the dispatcher in a FIFO order, and the *EUs* of their graph are ordered in breadth first order, while respecting their synchronization attributes (i.e., resource, condition and precedence). An *EU* eu_i is *runnable* if it meets the following four eligibility conditions: (i) the *EUs* that eu_i must wait for, due to precedence constraints, have finished their execution, (ii) all the resources needed by eu_i can be granted to eu_i , (iii) all the condition variables that eu_i must wait for

are set, (iv) the current time is higher than eu_i earliest start time. Among runnable *EUs*, the dispatcher executes the one with the highest priority.

The monitoring functionality of the dispatcher consists in checking tasks execution progress in order to detect the violation of safety, liveness and timeliness properties (e.g., missed deadline, violation of the task arrival law).

Kernel adaption layer: To port the HADES run-time support on different real-time kernel, a *kernel adaption layer* has been designed. This layer interacts with the underlying real-time kernel for activities related to threads and interrupts management. The adaptability of this kernel adaption layer is shown in Section 3.3.4.

Services layer: The services layer consists in a panel of services exhibiting fundamental properties (i.e., availability, dependability, timeliness) meeting the requirements of distributed dependable real-time applications. A *service* is defined as a module accomplishing a given function, and is called through a well-defined *interface*. A service interface is made of *task invocations*, and *upcalls* (i.e., task invocation requested by the service itself).

Services pertaining to the service layer are relevant to scheduling policies, distributed execution control, message multicast, group membership, time management, and synchronization of clocks (see Figure 1). In contrast to the dispatcher functionalities, these services are extensible, i.e., they can be modified/updated by the application designer as long as their implementation conforms to the service interface (see Section 3.2).

3.2 Support for extensibility

Extensibility is the property of systems in which the evolution of existing functionalities or the design of new ones is possible at any time during the life of the system. In order to achieve it, the dispatcher is decoupled from the services layer. This enables the application designer to extend the properties of any given service without having to cope with timing, synchronization, fault-tolerance and distribution constraints that are inherent to any other services. The following sections present the interface of four services, namely the schedu-

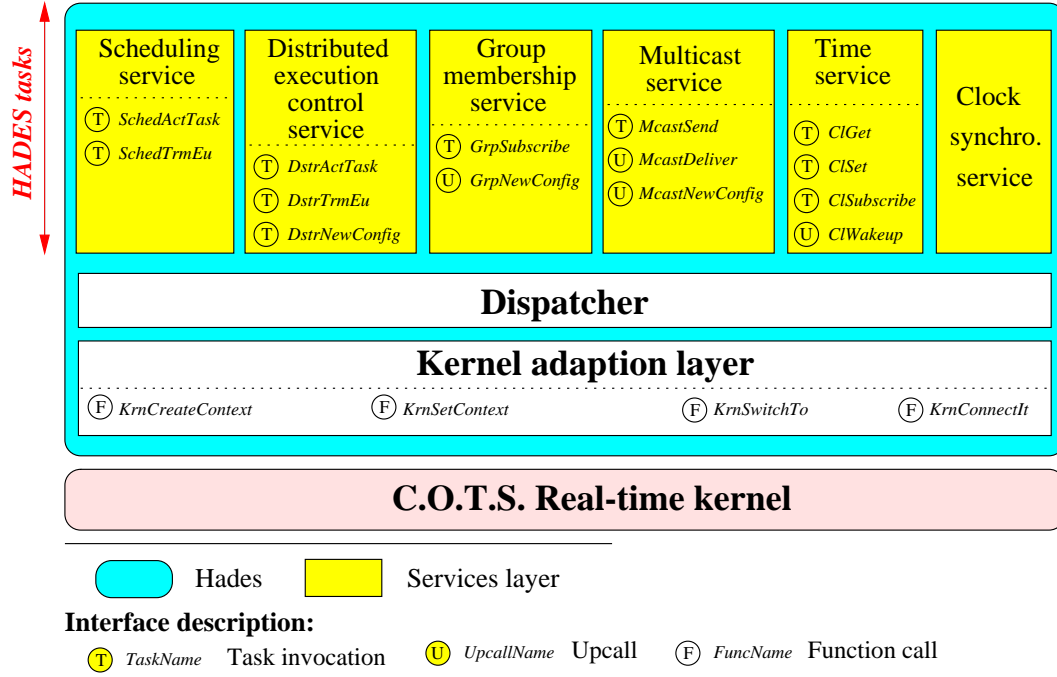


Figure 1: Internal structure of HADES run-time support

ling service, the distributed execution control service, the multicast service, and the group membership service.

Scheduling service: The scheduling service is in charge of computing the order in which *App_tasks* must execute to meet their performance requirements (e.g., deadline or priority). The feasibility tests, analyzing whether a set of tasks can meet its timing constraints or not, are provided by an off-line tool. In order to bind a given scheduling policy with a given *App_task* T_i , one indicates in the timing attribute of T_i 's *HEUG* (see Section 2) which scheduling policy to use. The scheduling service interface is composed of two *RTS_tasks* invocations, named *SchedActTask* and *SchedTrmEu*:

- Upon activation of an *App_task*, say T_i , the dispatcher asynchronously invokes the *SchedActTask* task, with two parameters identifying the task to be scheduled (i.e., T_i) and the set of T_i 's *EUs* that must be executed on the current site.
- Upon termination of an *EU*, say eu_i , the dispatcher asynchronously invokes the *SchedTrmEu* task with eu_i as argument.

The *SchedActTask* and *SchedTrmEu* tasks interact with the dispatcher via a set of *Sys_EUs* that manages the priority-ordered set of run-queues (i.e., insert first, insert last, insert after, and delete *Sys_EUs*).

The scheduling service interface has been designed to provide different scheduling policies, each being adapted to a different application domain (see Section 3.3.1).

Distributed execution control service: The distributed execution control service is in charge of launching the execution of distributed tasks on the relevant sites, validating their precedence constraints and detecting their termination. The execution control service is called by the asynchronous invocations of three *RTS_tasks*, named *DstrActTask*, *DstrTrmEu* and *DstrNewConfig*. The first two are in charge of respectively launching the distributed execution of *App_tasks*, and terminating their execution. The third one acknowledges modifications of the group's membership of the system (see hereafter). The distributed execution

control service to be used by a task T_i is specified in the distribution attribute of T_i 's *HEUG* (see Section 2). The distributed execution control service exhibits the following interface:

- Upon activation of a distributed *App_task* T_i , the dispatcher asynchronously invokes the *DstrActTask* task;
- Upon termination of an *EU* eu_i , the dispatcher asynchronously invokes the *DstrTrmEu* task;
- Upon modification of the group membership (i.e., after a site failure or site reinsertion), the *DstrNewConfig* task is asynchronously invoked by the group membership service, with two parameters indicating the new membership group view, and the time at which reinserted sites are bound to be operational. Notice that the *DstrNewConfig* task needs to be invoked only when the distributed execution control service has to meet a fault tolerant behavior.

The execution control tasks interact with the dispatcher via a set of *Sys_EUs* designed to interact with the core of the run-time support (e.g., allocation of an execution flow, validation of a set of precedence constraints, completion of the execution of an *invocation Sys_EU*). These tasks also interact with the multicast service to launch the task executions on the appropriate sites.

The execution control service interface has been designed to allow the implementation of execution control policies assuming a reliable and a fault-tolerant environment. For space consideration, these implementations are not given in this paper, but the interested reader is invited to read [CPC98].

Multicast service: The HADES multicast service is concerned with the transmission of messages to a set of sites. The HADES multicast service interface is made of two *RTS_tasks* invocations, named *McastSend* and *McastNewConfig*, and one *RTS_task* upcall named *McastDeliver*. The *McastSend* task broadcasts the messages to the sites belonging to the group while the *McastDeliver* upcall delivers the messages destined to the site. More precisely, the *McastDeliver* upcall invokes the appropriate tasks in charge of handling the

received messages. The *MCastNewConfig* task plays the same role as for the distributed execution control service. The multicast service interface is as follows:

- Upon request to multicast a message, the requesting task asynchronously invokes the *McastSend* task with three arguments indicating the message body, the destination sites and the task in charge of handling the message on each destination site.
- Upon receipt of a destined message, the *McastDeliver* upcall invokes the appropriate task in charge of handling the received message, with the message as an argument. The invocation of this task can be delayed to ensure that the delivery ordering properties of the multicast service are met (i.e., causal order, total, or temporal order).
- Upon modification of the group membership, the *McastNewConfig* task is asynchronously invoked with two arguments indicating the new membership group view and the time at which reinserted sites are bound to be operational.

The multicast service interface has been designed to allow the implementation of multicast protocols with different reliability and delivery order properties (see Section 3.3.2).

Group membership service: The HADES group membership service is in charge of detecting failures, allowing operational sites to agree on the set of failed and operational sites. Its interface is composed of a *RTS_task* invocation named *GrpSubscribe* allowing a HADES service to be notify upon group membership changes, and a *RTS_task* upcall named *GrpNewConfig* allowing to warn the subscribed services of changes in the group membership. More precisely:

- Upon request to be notify of group membership modifications, the requesting service invokes the *GrpSubscribe* task by indicating its identifier;
- Upon modification of the group membership (i.e., following a site failure, or site reinsertion), the *GrpNewConfig* upcall provides the new membership group view as well as the time at which reinserted sites are bound to be operational to subscriber services.

Notice that the respective interfaces of the multicast and group membership services do not impose the role and number of tasks used to implement these services. For performance reasons, the service designer can implement these two services jointly by using the same task to implement both the multicast and group membership protocols (see section 3.3.2 and 3.3.3), while preserving their flexibility.

3.3 Support for adaptability

HADES specifies the interface of a set of run-time support services, so that the properties they exhibit can be adapted to the target application domain, to the hardware, or to the fault hypothesis, as long as they conform to the service interface (see Section 3.2). This section illustrates this flexibility by giving for three services of the Services layer, two different specifications according to the scheduling requirements (Scheduling service), the message delivery order (Multicast service), the group membership properties (Group membership service), and the underlying kernel features (Kernel adaption layer).

3.3.1 Scheduling service

The choice of a scheduling service highly relies on application-specific features like task arrival laws, task priorities, preemption policy, and resource access patterns. Several dynamic and static scheduling algorithms (e.g. Earliest Deadline First (EDF), Rate Monotonic (RM) [LL73]) as well as planning-based scheduling policies [XP90, RSS90, Agn91]) exist and some of them have been specified and implemented in the HADES run-time support. Two scheduling policies that conform to the interface of the scheduling service are presented hereafter: a static *planning-based scheduling* policy and a *dynamic scheduling* policy. The first implementation is intended for applications whose required resources must be preallocated so that deadlines can be guaranteed *a priori*. The motivation for the second implementation is to support sporadic tasks.

Static planning-based scheduling policy: The static planning-based scheduling policy we have specified and implemented is based on the Xu and Parnas algorithm [XP90], which

acts on a set of periodic tasks with resources and precedence constraints. The original algorithm has been extended to handle distributed tasks. This scheduling policy relies on a feasibility test run off-line, which generates a scheduling plan, containing for every *EU* its start time relative to the beginning of the task period. A set of HADES run-queues is reserved for *EUs* scheduled by the extended Xu and Parnas scheduler. Each *EU* is mapped onto one of this run-queue depending on its preemption level (i.e. if *EU* eu_i preempts *EU* eu_j , it is assigned to a run-queue with a higher priority). The execution of the Xu and Parnas scheduling policy is implemented by two *RTS_tasks* named $t_{XuParnas}^{Act}$ and $t_{XuParnas}^{Trm}$. Upon activation of an *App_task*, that is at the beginning of the task period, task $t_{XuParnas}^{Act}$ inserts each *EU* of this task into the dispatcher run-queues. It uses the information given in the scheduling plan to insert the *EUs* in the run-queue corresponding to their preemption level and for a given run-queue, to sort them according to their increasing earliest start times. Upon *EU* termination, task $t_{XuParnas}^{Trm}$ removes the *EU* from the run-queue it belongs to.

Dynamic scheduling policy: The dynamic scheduling policy we have implemented is the Earliest Deadline First (EDF) scheduling policy coupled with the Stack Resource Policy (SRP) [Bak91] to avoid multiple priority inversions. The key idea of SRP is that when a task needs a resource, it is blocked at start time rather than later when it actually requests the shared resource. The main motivation for this early blocking is to save unnecessary context switches. Briefly, the SRP statically assigns a *preemption level*, noted $\pi(t)$ to each task t . Preemption levels are defined so that they are ordered inversely with respect to the order of relative deadlines.

The implementation of the EDF scheduling policy without the SRP is straightforward. The *App_tasks* scheduled according to the EDF policy are assigned to a single HADES run-queue, into which tasks are sorted by increasing deadlines. The EDF scheduling policy is managed by two *RTS_tasks* named t_{edf}^{Act} and t_{edf}^{Trm} . Each time a task is activated, task t_{edf}^{Act} adds the *EUs* of this task to the EDF run-queue according to their deadlines. Each time an *EU* terminates, task t_{edf}^{Trm} removes it from the EDF run-queue. This implementation undergoes some modifications when the EDF scheduling policy is coupled with the SRP.

The EDF/SRP scheduling policy is managed by two *RTS_tasks* named $t_{edf/srp}^{Act}$ and $t_{edf/srp}^{Trm}$. Each time a task t is activated, task $t_{edf/srp}^{Act}$ checks if t can preempt the currently executing task t_{cur} . According to SRP, t can preempt t_{cur} if the following preemption condition is met: $deadline(t) < deadline(t_{cur})$ and $\bar{\pi} < \pi(t)$, where $\bar{\pi}$ is the maximum of the current ceilings of all resources [Bak91]. If the preemption condition is met, task $t_{edf/srp}^{Act}$ inserts the *EUs* belonging to t at the beginning of the EDF/SRP run-queue; otherwise, t is kept in an internal data structure of $t_{edf/srp}^{Act}$, containing pending task execution requests. Each time an *EU* terminates, task $t_{edf/srp}^{Trm}$ removes it from the EDF/SRP run-queue and scans its list of pending task execution requests to test the preemption condition and to modify the EDF/SRP run-queue if necessary.

3.3.2 Time-bounded reliable multicast

The problem of sending messages in a failure-prone environment has been subject to intensive investigation, yet relatively little research seems to address the problem of designing reliable multicast that are amenable to schedulability analysis. HADES run-time support provides three time-bounded multicast protocols. The first one, called *time-bounded basic multicast* protocol, has been designed for a failure-free environment, whereas the others rely on the fault hypothesis given in Section 2.3. These two fault tolerant multicast protocols, named *time-bounded atomic multicast* protocol and *time-bounded causal multicast* protocol ensure that any message sent by a correct site will be delivered by all correct destination sites by guaranteeing the *Validity*, *Integrity*, *Agreement* and *real-time Δ -Timeliness* properties [HT93]. They only differ by the order in which messages are delivered: *temporal order* for the first one and *causal order* for the second protocol. The causal order ensures that messages are delivered in accordance with their precedence relationship. The temporal order of events in a distributed system is defined in terms of an omniscient external observer. This observer possesses a single reference clock and timestamps each event by this clock. So event E_1 temporally precedes event E_2 if $timestamp(E_1) < timestamp(E_2)$. In simpler terms, if E_1 occurs before E_2 , then E_1 precedes E_2 .

For space considerations, we only give a brief description of the atomic and the causal time-bounded multicast protocols.

Time-bounded atomic multicast protocol: The HADES real-time atomic multicast protocol is based on the Time-Division Multiple Access strategy, and relies on the clock synchronization service to get local clocks synchronized within γ of each other. Time is divided into rounds, a round being composed of a succession of time slots. Each site i gets s_i slots per round to send its messages. The assignment of slots to sites is done statically; it is never modified even in case of group membership changes. Upon message submission, task *McastSend* stores the message in a sending buffer. Concurrently, a periodic *RTS_task* called $t_{send/recv}$ is activated during the time slots granted to the site. This task is in charge of timestamping each message of the sending buffer with a unique identifier based on the local time of the site, broadcasting them on the channel, and receiving each broadcast message. Note that received messages are piggybacked in broadcast messages until each message has been broadcast $\omega = f_{send} + 1$ times to tolerate up to f_{send} successive omissions. For each destined message m , the $t_{send/recv}$ task invokes the task in charge of handling m (i.e., upcall *McastDeliver*) as soon as m 's timestamp satisfies the delivery conditions. These conditions concern the respect of the temporal order. This protocol ensures that if a message m is multicast at real-time t (by invoking task *McastSend*), then no correct site invokes the destination task in charge of handling m after real-time $t + \Delta$. The value of Δ is equal to $2((\omega + 2).(P + r + 2.\omega) + \Delta_{trans} + 2.\omega)$, with P the period of task $t_{send/recv}$, Δ_{trans} the maximum network latency in a failure-free environment, γ the precision of the clock synchronization algorithm, and r the maximum response time of task $t_{send/recv}$. By taking, $\Delta_{trans} = 1ms$, $P = 10ms$, $r = 1ms$, $\gamma = 10ms$, and $\omega = 2$, we get $\Delta \simeq 290ms$.

Time-bounded causal multicast protocol: The HADES real-time causal multicast protocol ensures that messages are causally delivered. Compared to the time-bounded atomic multicast protocol, this protocol does not rely on synchronized clock which offers an increased flexibility. This protocol is implemented by a single periodic *RTS_task* named $t_{send/recv}$. This task is periodically activated, and broadcasts ω times each message of the sending buf-

fer. Upon receipt of a message m , $t_{send/recv}$ invokes the task in charge of handling m when delivery conditions are satisfied. To avoid infinite messages collisions on the communication medium, a mechanism providing a unique access to the medium is mandatory (e.g., in HADES, the ATM switch is used; in case where Ethernet were used instead of ATM, the medium access protocol CSMA/CD would guarantee a deterministic access [HLR95]). The value of Δ is for this protocol equal to $(3.\omega + 1).P + \Delta_{trans} + 2.r$. By taking the above numerical values, we get $\Delta \simeq 73ms$.

3.3.3 Group membership service

The challenge when designing a group membership protocol for a hard real-time system is to fulfill the three following requirements: (i) to have a small and bounded failure detection latency, (ii) to have a small and bounded number of messages exchanged and (iii) to have a protocol that can be easily integrated into feasibility tests. The HADES run-time support comes with two different implementations of the group membership service that meet such constraints. The first implementation, called *strong* group membership, guarantees the *Vivacity*, *Uniqueness*, *Consistency* and *Real-time Δ -Timeliness* properties [Cri91]. For performance reasons, the second one, called *weak* group membership service, guarantees only the *Vivacity*, and *Real-time Δ -Timeliness* properties. Both implementations periodically piggybacks heartbeat messages on application messages (in the current implementation, the $t_{send/recv}$ task is used to implement both the multicast and the group membership services). For space consideration, the following is devoted to the description of only the strong group membership protocol.

The group membership service is based on the same principle as the time-bounded atomic multicast service, i.e., a TDMA approach. Each site p piggybacks on each application message the view it has of the sites group, by means of two vectors named $recv_p$ and $join_p$. Vector $recv_p$ is used to implement the heartbeat: $recv_p[i]$ is the latest clock value that p knows that site i sent. Vector $join_p$ is used to agree on when a site has restarted: when the current local time is at least $join_p[i]$, then i has joined the group of operational sites. Upon receipt of a message on site p , vector $recv_p$ is updated and the message is forwarded to the

group. This protocol guarantees that if at local time t_p , a site p operational for more than Δ time units that has not received any heartbeat message from site q for Δ time units, then p excludes q from its group view at local time t_p . The value of Δ is equal to the one exhibited by the time-bounded atomic multicast protocol (see Section 3.3.2). Thanks to the clock synchronization service, all correct sites exclude q from their group view at time $t_q = t_p$. Similarly, a site p initially crashed, and wishing to be reinserted in the group is guaranteed to be reinserted in the new group view by all correct sites at time $t_p + \Delta$, with t_p the time at which p broadcast its first heartbeat message ($join_p[p] = t_p + \Delta$). Upon configuration change, tasks that need to be notified of the configuration are invoked with the new group view as argument. Note that the time-bounded atomic multicast service altogether with the group membership service ensure the virtual synchrony property, i.e., the guarantee that the delivery of messages are totally ordered with respect to changes in the group's membership.

3.3.4 Kernel adaption layer

The kernel adaption layer interacts with the underlying real-time kernel for activities related to the management of threads and interrupts. To port the HADES run-time support on a real-time kernel, one has only to develop a kernel adaption layer that conforms to a interface made of four function-calls directly invoked by the dispatcher:

Name	Description
KrnCreateContext	Allocates a context for a thread
KrnSetContext	Sets the thread context (e.g., initialize the thread starting address) but does not start the thread yet.
KrnSwitchTo	Starts a thread or restarts it after a preemption
KrnConnectIt	Connects an interrupt handler on a given interrupt level.

The first three function calls allow the dispatcher to assign separate threads of control to concurrent *EUs* within a *HEUG*. Upon occurrence of an interrupt, the fourth primitive connects the interrupt handler of the concerned HADES services (e.g., time service).

Such an interface requires that the real-time kernel offers enough support *(i)* to implement a multithreading interface; *(ii)* to connect user-defined interrupt handlers. In addition the real-time kernel must guarantee the consistency of the dispatcher in presence of interrupts, and this without any blocking. This last point is detailed hereafter. Up to now, the HADES kernel adaption service has been ported to Chorus 3.1 and has been emulated on Solaris.

Chorus kernel adaption service: Chorus 3.1 [CS96] provides a real-time operating system that can scale down to small embedded platforms and scale up to distributed POSIX compliant platforms. It is built around a *core executive* module which provides the basic services mandatory for a real-time executive. The HADES environment has been ported on the *micro core* executive module. This module provides support for multithreaded applications, FIFO fixed-priority scheduling, basic synchronization and user-provided management of interrupts, traps and exceptions. Note that the *core executive* module, while offering more functionalities (e.g., memory protection, dynamic memory management) has not been selected for predictability reasons (its code was far too large to be guaranteed predictable through manual code analysis). The HADES kernel adaption service interface has directly been implemented above the *micro core* API.

The most difficult part concerned the implementation of indivisible portions of code on top of the *micro core*. One trivial solution could have been to simply mask interrupts (using *ix86* instructions *cli* and *sti*) during the portions of code that had to be indivisible. This solution was rejected since interrupts are unmasked at the end of all Chorus system calls, and thus prevent portions of code invoking them to be indivisible. Our solution relies on the *micro core* fixed-priority scheduler. It consists in assigning Chorus priorities to thread execution according to the following policy: *(i)* threads allocated to *Sys_EUs* are executed with the highest Chorus priority; *(ii)* threads allocated to interrupt handlers are executed with a Chorus intermediate priority. These threads are unblocked whenever an interrupt is triggered; *(iii)* threads allocated to *Code_EU* are executed with a low Chorus priority. With such a priority assignment, an interrupt can be triggered while a *Sys_EU* is being executed, but its interrupt handler is executed only when the *Sys_EU* is terminated. This solution

can be implemented on top of any real-time kernel exhibiting a fixed-priority scheduling policy.

Solaris kernel adaption service: The HADES kernel adaption service has been ported on the Solaris operating system. This port aimed at validating HADES temporal properties by simulating a worst-case execution scenario, and HADES functional properties. The HADES kernel adaption layer interface has directly been implemented above the Solaris thread library. More precisely, per each site is allocated one Solaris multi-thread process in charge of executing all the *App_tasks* and *RTS_tasks*. For the global system is allocated one Solaris multi-thread process in charge of (i) coordinating all the sites, (ii) simulating the global time, (iii) handling the local clocks drift, and the messages transmission time.

3.4 Support for predictability

A key design principle in HADES is to provide a run-time system with a *predictable* behavior. An action is said to be predictable if its results and its worst case execution time (WCET) are known before it is executed. Predictability is crucial for the implementation of feasibility tests. Moreover, the more precise the WCET estimation is, the less pessimistic the application feasibility test is. Indeed, due to the complexity of determining worst case cost information, scheduling tests often use over-estimated WCETs of run-time support activities. While this behavior is safe, it often leads to a negative answer of the scheduling test, forbidding the execution of applications despite their actual feasibility. Providing a predictable run-time support requires that all its components be predictable. A WCET analyzer is currently being developed to partially automate kernel and *Code_EUs* WCET analysis. The following sections show how the design of all components of the HADES environment (i.e., services layer, dispatcher and real-time kernel) makes HADES predictable.

Services layer predictability The predictable behaviour of the service layer results from two aspects. The first one comes from the specification of all run-time support services according to the HADES task model. This model eases the determination of the temporal behavior of these services by requiring that sequential codes (i.e., *EUs*) are isolated from

any synchronization, or distribution constraints. This makes possible a precise evaluation of *EUs WCET*, and thus an exact knowledge of any blocking time. The second aspect is related to the arrival law of these services. Except for a small set of *RTS_tasks* which are periodical (the $t_{send/recv}$ task of the multicast service, the clock synchronization and the time tasks), all the services are invoked explicitly or implicitly by the application code. This enables to carry their WCETs over to the execution cost of application tasks. The periodic laws of the aforementioned *RTS_tasks* are directly integrated into the scheduling feasibility test.

Dispatcher predictability The dispatcher has been designed to handle very basic functionalities, which make its code small and easy to analyze. Beyond this feature, the dispatcher exhibits a number of characteristics that make it predictable. More precisely, the dispatcher (i) does not contain any asynchronous or periodic activity; (ii) does not use dynamic memory allocation (tasks arrival laws, resources, condition variables are known off-line); and (iii) never blocks (see Section 3.3.4).

Such features enable to fully characterize the dispatcher with four WCETs: (i) $wcet_{term/start}$ to terminate an *EU* execution and to start the execution of the following one; (ii) $wcet_{sys}$ to execute a *Sys_EU*; (iii) $wcet_{int}$ to take into account an interrupt and to start the execution of the appropriate interrupt handler; (iv) $wcet_{except}$ to substitute the execution of the current *Code_EU* by the exception handler. These four WCETs have been identified in our prototype and correspond to treatments that cannot be preempted. The integration of these costs into the application feasibility test is relatively easy. The $wcet_{term/start}$ and $wcet_{sys}$ costs are directly added to the application tasks costs, while the $wcet_{int}$ and $wcet_{except}$ are integrated as high priority sporadic tasks.

Kernel predictability As the designers of real-time kernels generally do not provide any guarantee about predictability of their kernel, access to their source code is required. Thus, a study of the Chorus R3 kernel adopted in the HADES prototype, has been done. The ChorusR3 kernel source code has been analyzed through a verification of invariants, an

analysis of function side effects, and an identification of WCETs. So far, code analysis has been done manually.

4 Related work

A number of systems have been designed to support dependable real-time distributed applications. Mars [K⁺89, RSL95] and Maft are two of them. However, Mars is not flexible: first it uses a single scheduling policy (static and periodical), and second it relies on specific hardware-intensive solutions for fault-tolerance and clock synchronization. In contrast, HADES run-time support can be adapted to support multiple scheduling policies (see Section 3.3) and uses COTS hardware and software. Maft [KWFT88] proposes mechanisms to tolerate byzantine failures in order to provide extremely reliable computations, and this, without sacrificing performances. However, for the same reasons as Mars, Maft is not flexible. Delta-4/XPA [BHV⁺90] proposes flexible fault-tolerance and communication protocols. Unfortunately, the execution environment used for its implementation is not predictable and the proposed fault-tolerance real-time protocols are not suited for hard real-time constraints. Among works in progress, the GUARDS European ESPRIT project [WBDBP96] has objectives similar to those of HADES. However, in opposition to HADES, GUARDS favors the development of hardware components easing the implementation of fault-tolerance mechanisms.

HADES run-time support cannot be directly compared with real-time kernels, like QNX [Hil92] and Chorus [CS96]. While some of them exhibit a predictable behavior, they are generally not designed to support failures or to adapt to application needs (e.g., support multiple scheduling policies). However, they can be used as HADES base real-time kernel so long as they offer enough functionalities to implement the HADES kernel adaption layer.

The structuring principles that were adopted in HADES to meet the *flexibility* property are very similar to the ones proposed by micro kernel designers to specialize their kernels [Lie96]. The difference comes mainly from the class of service that can be specialized. In micro kernels, services that can be specialized are general-purpose operating system services (e.g., memory management specialization through paging services, communication services), while in HADES, such services are representative of services found in run-time

supports for dependable applications (e.g., real-time scheduling, time-bounded communications, clock synchronization). More recent work on kernel specialization gave birth to highly specializable operating systems kernels, such as Spin [B⁺95] or Exokernel [K⁺97], but these kernels do not address predictability issues.

CORBA [OMG96] is a standardized middleware architecture for distributed object computing on heterogeneous environments. It eases the development of distributed services by providing features to interconnect applications and services. TAO [SLM97] is a predictable implementation of CORBA, providing facilities to support hard real-time requirements. HADES, like TAO relies on a predictable middleware layer. However the two platforms differ by the fact that HADES does not rely on an object-based approach to develop applications, by the fact that HADES does conform to any standard like CORBA and by the fact that to our knowledge, TAO is not fault-tolerant.

References

- [ACCP98] E. Anceaume, G. Cabillic, P. Chevochot, and I. Puaut. Hades: A middleware support for distributed safety-critical real-time applications. In *Proc. of ICDCS-18*, pages 344–351, May 1998.
- [Agn91] R. Agne. Global cyclic scheduling : A method to guarantee the timing behavior of distributed real-time systems. *The Journal of Real-Time Systems*, 3(1):45–66, March 1991.
- [B⁺95] B. Bershad et al. Extensibility, safety and performance in the SPIN operating system. In *Proc. of SOSP-15*, volume 29(5) of *OSR*, pages 267–284, October 1995.
- [Bak91] T. Baker. Stack-based scheduling of real-time processes. *The Journal of Real-Time Systems*, 3(1):67–99, 1991.

- [BHV⁺90] P. Barrett, A. Hilborne, P. Verissimo, L. Rodrigues, P. Bond, D. Seaton, and N. Speirs. The delta-4 extra performance architecture (xpa). In *Proc. of FTCS-20*, pages 481–488, June 1990.
- [CPC98] P. Chevochot, I. Puaut, and G. Cabillic. A flexible environment for automatic replication and execution of distributed safety-critical real-time applications. TR 1183, IRISA, April 1998.
- [Cri91] F. Cristian. Reaching agreement on processor group membership in synchronous distributed systems. *Distributed Computing*, 4(4):175–187, 1991.
- [CS96] Chorus Systèmes. Chorus/classix R3 technical overview. TR 96-119.91, May 1996.
- [Hil92] Dan Hildebrand. An architectural overview of QNX. In *Proc. of the USENIX Workshop on Micro-Kernels and Other Kernel Architectures*, pages 113–126, April 1992.
- [HLR95] J.F. Hermant, G. Le Lann, and N. Rivierre. A general approach to real-time messages scheduling over distributed broadcast channels. In *IEEE Conf. on Emerging Technologies and Factory Automation*, October 1995.
- [HT93] V. Hadzilacos and S. Toueg. *Distributed systems, fault tolerant broadcasts and related problems*. Addison-Wesley / ACM Press, 1993.
- [K⁺89] H. Kopetz et al. Distributed fault-tolerant real-time systems: The Mars approach. *IEEE Micro*, 9:25–40, February 1989.
- [K⁺97] M. Kaashoek et al. Application performance and flexibility on exokernel systems. In *Proc. of SOSP-16*, volume 31(5) of *OSR*, pages 52–65, October 1997.
- [KWFT88] R. Kieckhafer, C. Walter, A. Finn, and P. Thambidurai. The MAFT architecture for distributed fault tolerance. *IEEE Trans. on Computers*, 37(4):398–405, April 1988.

- [Lie96] J. Liedtke. Toward real microkernels. *Comm. of the ACM*, 39(9):70–77, September 1996.
- [LL73] C.L. Liu and J. Layland. Scheduling algorithms for multiprogramming in a hard real-time environment. *Journal of the ACM*, 20(1):46–61, January 1973.
- [MRS⁺90] L. Molesky, K. Ramamritham, C. Shen, J. Stankovic, and G. Zlokapa. Implementing a predictable real-time multiprocessor kernel - the Spring kernel. In *Proc. of the 7th IEEE Workshop on Real-Time Operating Systems and Software*, May 1990.
- [OMG96] OMG. *The Common Object Request Broker: Architecture and Specification*, July 1996.
- [RSL95] J. Reisinger, A. Steininger, and G. Leber. *Predictably Dependable Computing Systems*, chapter The PDCS Implementation of MARS Hardware and Software, pages 209–224. Springer-Verlag, 1995.
- [RSS90] K. Ramamritham, J. Stankovic, and P. Shiah. Efficient scheduling algorithms for real-time multiprocessor systems. *IEEE Trans. on Parallel and Dist. Systems*, 1(2):184–194, April 1990.
- [SLM97] D.C. Schmidt, D.L. Levine, and S. Mungee. The design of the TAO real-time object request broker. *Computer Comm. Journal*, 1997.
- [WBDBP96] A. Wellings, L. Beus-Dukic, A. Burns, and D. Powell. Genericity and upgradability in ultra-dependable real-time architectures. TR LAAS-R-96417, LAAS, November 1996.
- [XP90] J. Xu and D. L. Parnas. Scheduling processes with release times, deadlines, precedence and exclusion relations. *IEEE Trans. on Software Engineering*, 16(3):360–369, March 1990.



Unité de recherche INRIA Lorraine, Technopôle de Nancy-Brabois, Campus scientifique,
615 rue du Jardin Botanique, BP 101, 54600 VILLERS LÈS NANCY

Unité de recherche INRIA Rennes, Irisa, Campus universitaire de Beaulieu, 35042 RENNES Cedex

Unité de recherche INRIA Rhône-Alpes, 655, avenue de l'Europe, 38330 MONTBONNOT ST MARTIN

Unité de recherche INRIA Rocquencourt, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex

Unité de recherche INRIA Sophia-Antipolis, 2004 route des Lucioles, BP 93, 06902 SOPHIA-ANTIPOLIS Cedex

Éditeur

INRIA, Domaine de Voluceau, Rocquencourt, BP 105, 78153 LE CHESNAY Cedex (France)

<http://www.inria.fr>

ISSN 0249-6399